

Data Science in Action – 10 Dinge, die Advanced Analytics und Data Science für Ihr Unternehmen tun kann

Gerhard Svolba, Analytic Solution Architect, SAS Austria
Wien, 11. Oktober 2017







Agenda

- 10 mal „Data Science in Action“
 - Supervised Machine Learning Methoden
 - Unsupervised Machine Learning Methoden
 - Simulationen
- Data Science und Advanced Analytics mit der SAS Analytic Plattform
- Zusammenfassung und Links

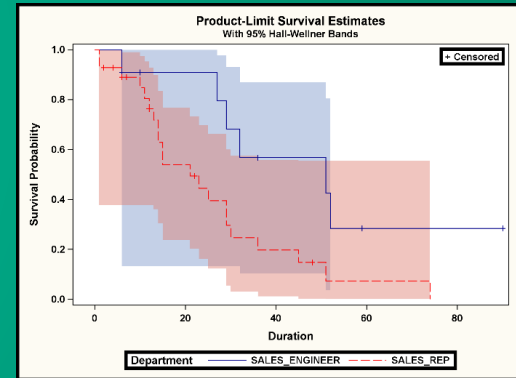
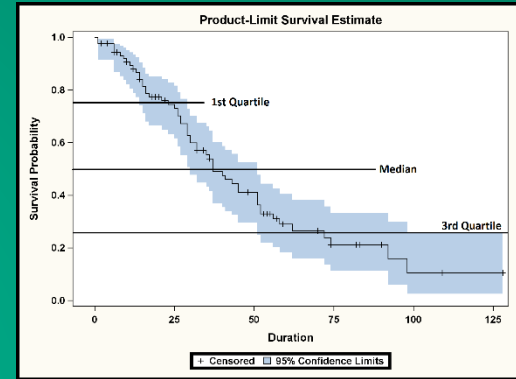


Data Science in Action: #1

Performing Headcount Survival Analysis for Employee Retention

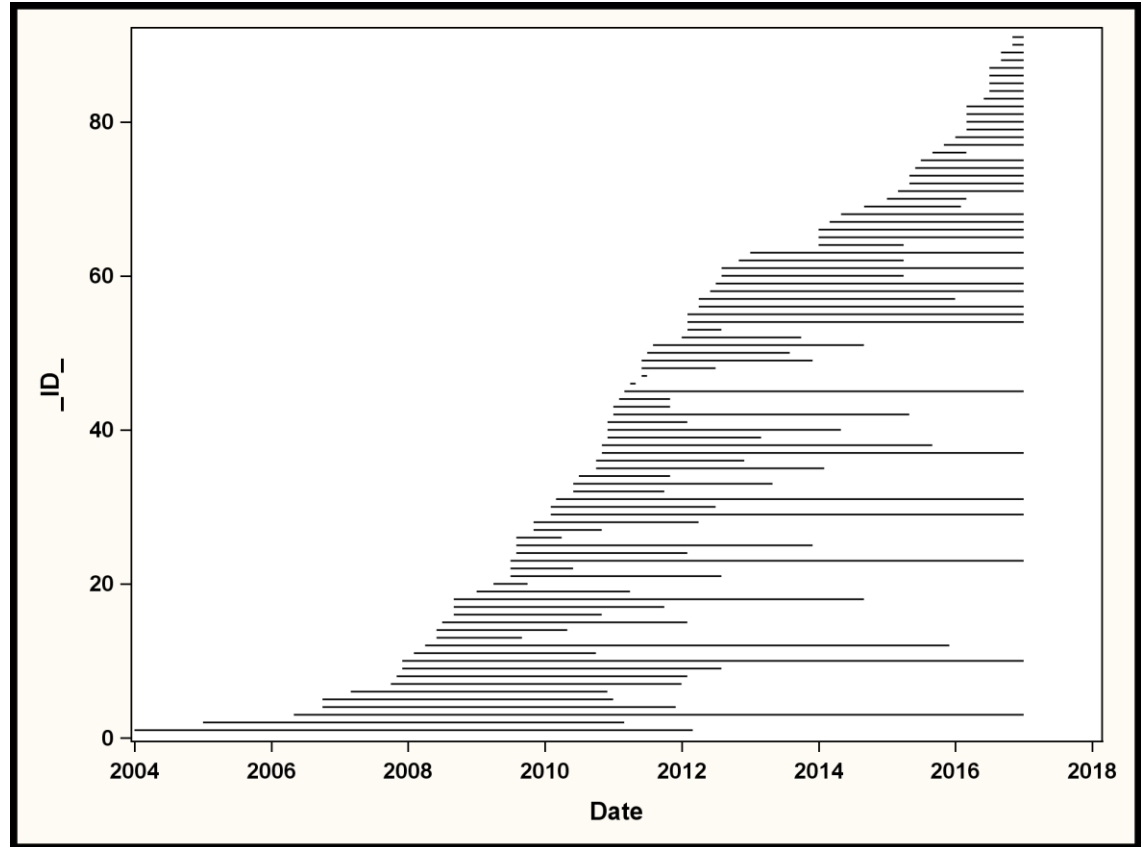
*Can assumptions about the average
length of time intervals be made, even if
most of the endpoints have not yet been
observed?*

Survival analysis methods: Kaplan-Meier estimates
Cox Proportional Hazards regression
Survival Data Mining

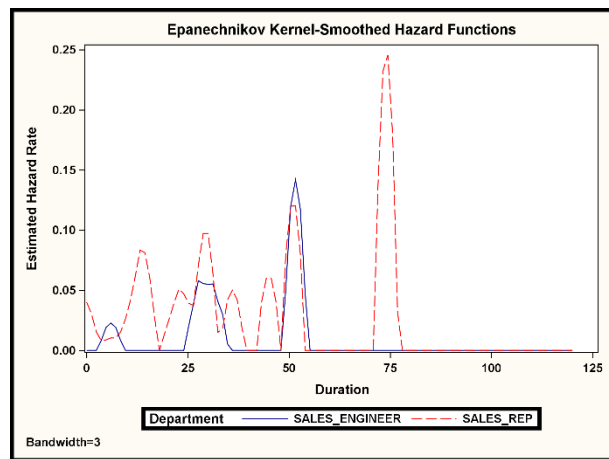
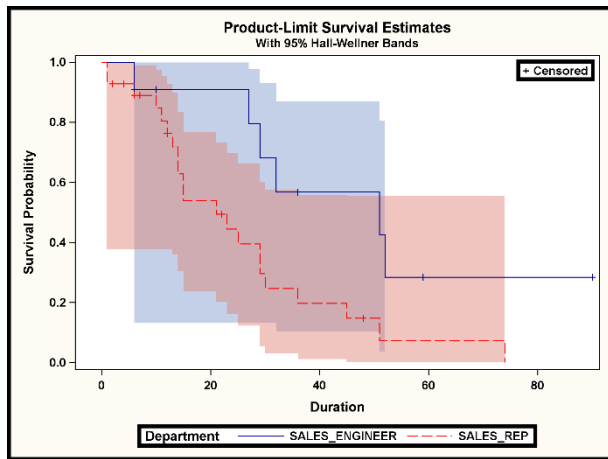


Nicht zu allen Mitarbeitern haben wir ein „Ereignis-Datum“ (Glücklicherweise)

- Betrachten der Karrieren pro Mitarbeiter
 - Unterschiedliche Länge
 - Kündigung oder „zensiert“



Die Kaplan Meier Methode und die Cox Proportional Hazards Regression verarbeitet zensierte Beobachtungen



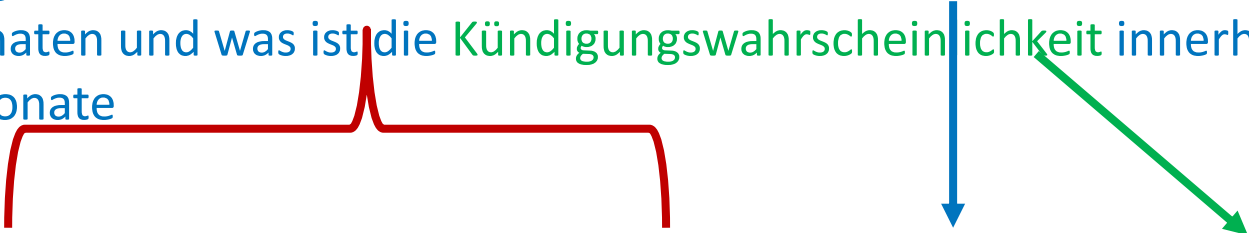
Kaplan Meier Methods und Cox Proportional Hazards Regression:
Sales engineers haben eine bessere „survival time“ als sales representatives.

Betrachten der Hazard Kurven:
Es gibt ein hohes Risiko die Sales Engineers nach 26 und 50 Monaten zu verlieren.

„Wie lange wird Gerhard Svolba noch in unserem Unternehmen sein?“

Vorhersage der Verweildauer für individuelle Mitarbeiter

Ausgehend von bestimmten Risikofaktoren, was ist die erwartete Survival in 6 Monaten und was ist die Kündigungswahrscheinlichkeit innerhalb der nächsten 6 Monate



EmpNo	Department	Gender	TechKnowH...	_T_	EM_SURVFCST	EM_SURVEVENT	T_FCST
1003	TECH_SUPPORT	M	YES	128	0.240	0.000	134
1010	TECH_SUPPORT	M	YES	109	0.240	0.011	115
1023	SALES_ENGINEER	M	YES	90	0.108	0.313	96
1029	TECH_SUPPORT	M	YES	83	0.386	0.133	89
1031	TECH_SUPPORT	F	YES	82	0.177	0.219	88
1037	ADMINISTRATION	M	NO	74	0.471	0.066	80
1045	ADMINISTRATION	M	NO	70	0.494	0.053	76
1054	TECH_SUPPORT	F	YES	59	0.316	0.102	65
1055	SALES_ENGINEER	M	YES	59	0.313	0.103	65

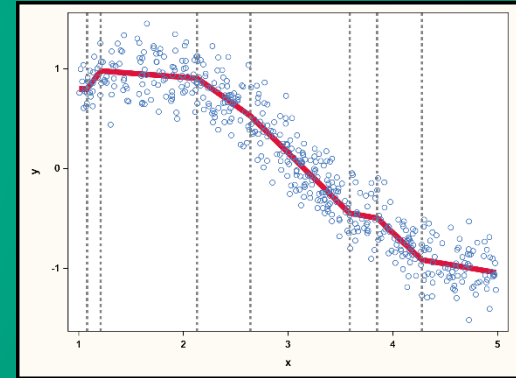
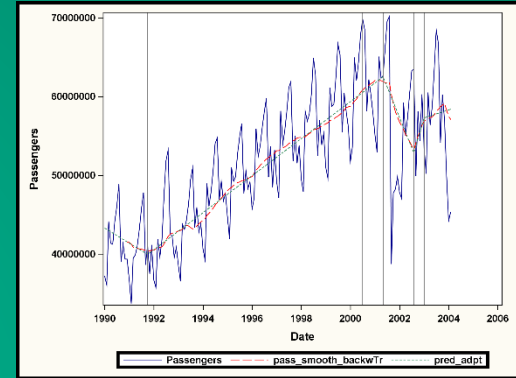


Data Science in Action: #2

Detecting Structural Changes and Outliers in Longitudinal Data

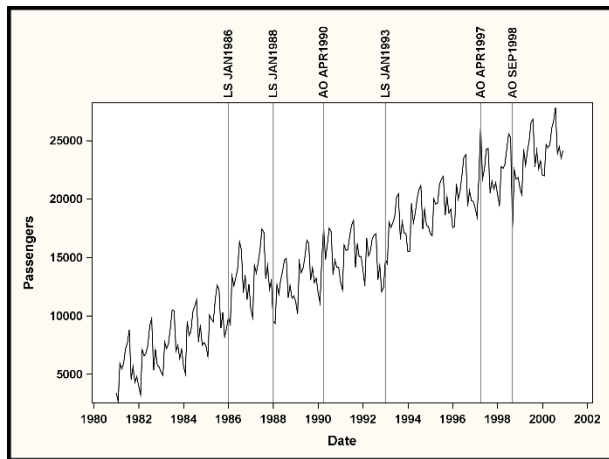
*Can events and changes in the
course over time be
automatically detected?*

Smoothing Of Longitudinal Data
Multivariate Adaptive Regression Splines
Automatic Breakpoint Detection
Automatic Detection of Outliers with ARIMA Models

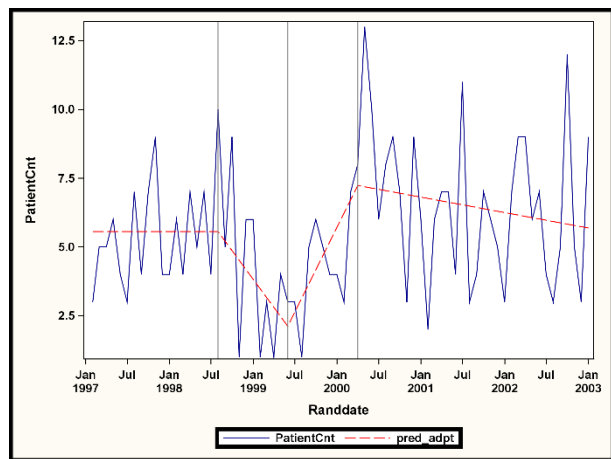


Automatisches Erkennen von Breakpoints und Ausreißern

Anwenden von analytischen Methoden zum Erkennen von Zeitpunkten, wo der Verlauf der Daten vom „normalen“ Muster abweicht.



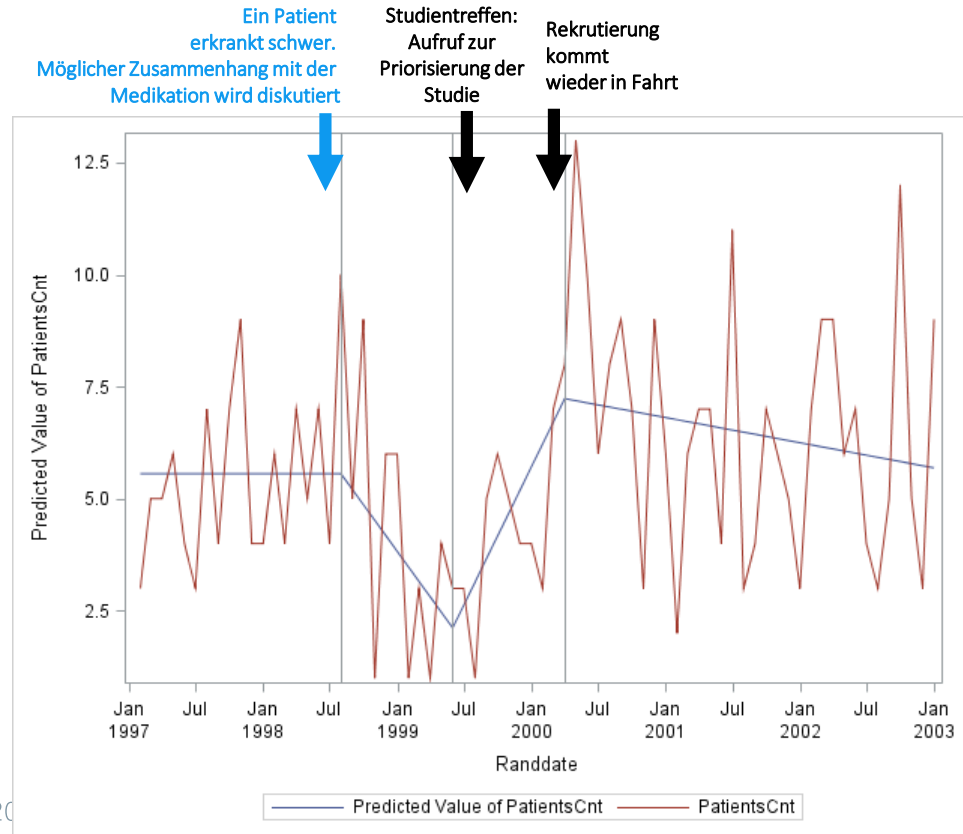
Erkennen von Shifts und Pulse Events mit ARIMA Modellen



Verwenden von Multivariate Adaptive Regression Splines zum Auffinden von Bruchpunkten



Was ist zu bestimmten Zeitpunkten in meiner klinischen Studie passiert?

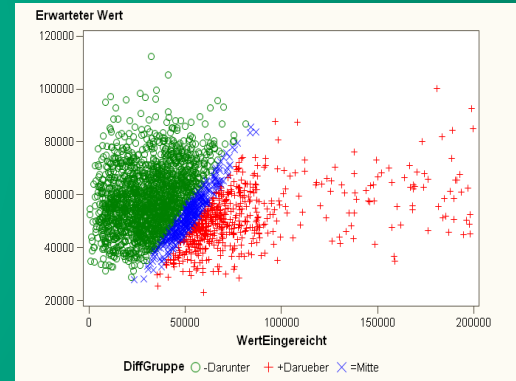
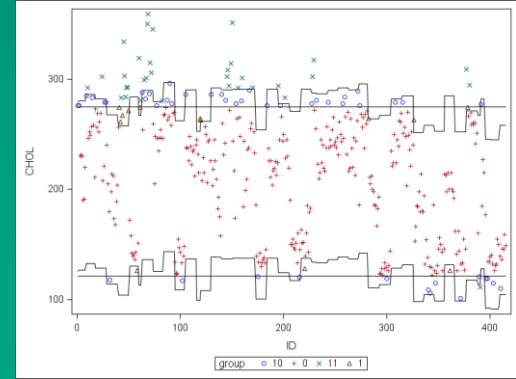


Data Science in Action: #3

Proving a reference value
that considers all available
co-information

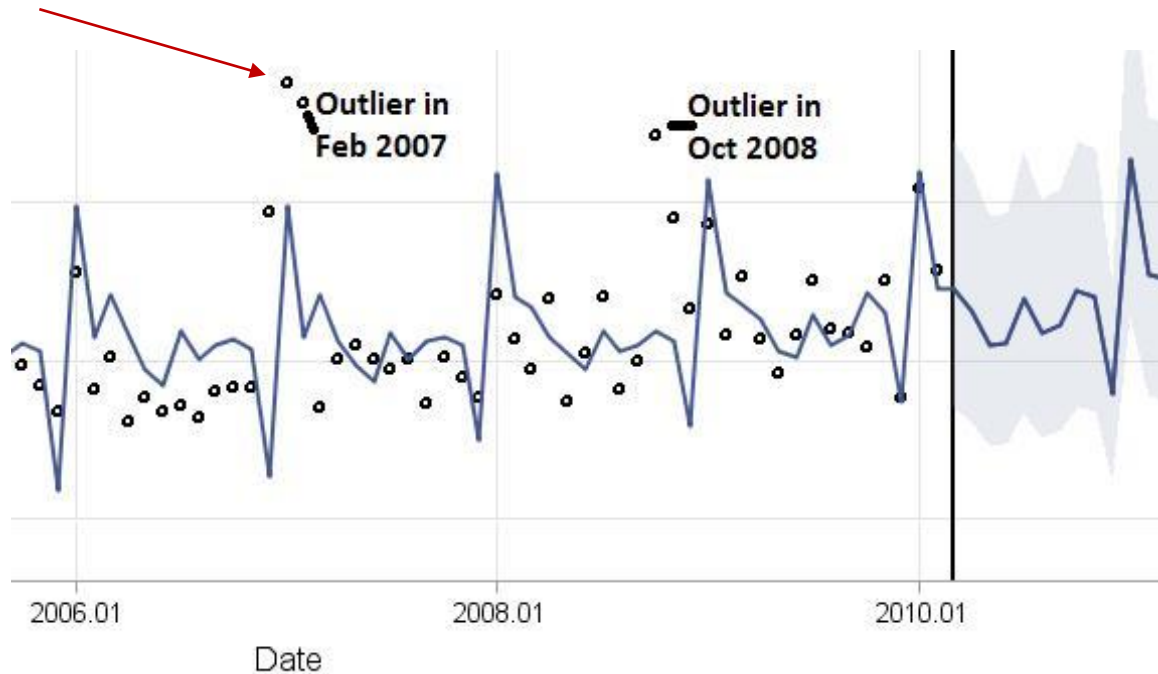
*Can analytics help me to reduce the
“Yes, but ... “ sentences in my business
dicussions?*

Linear Regression
Decision Trees
Time Series Analysis

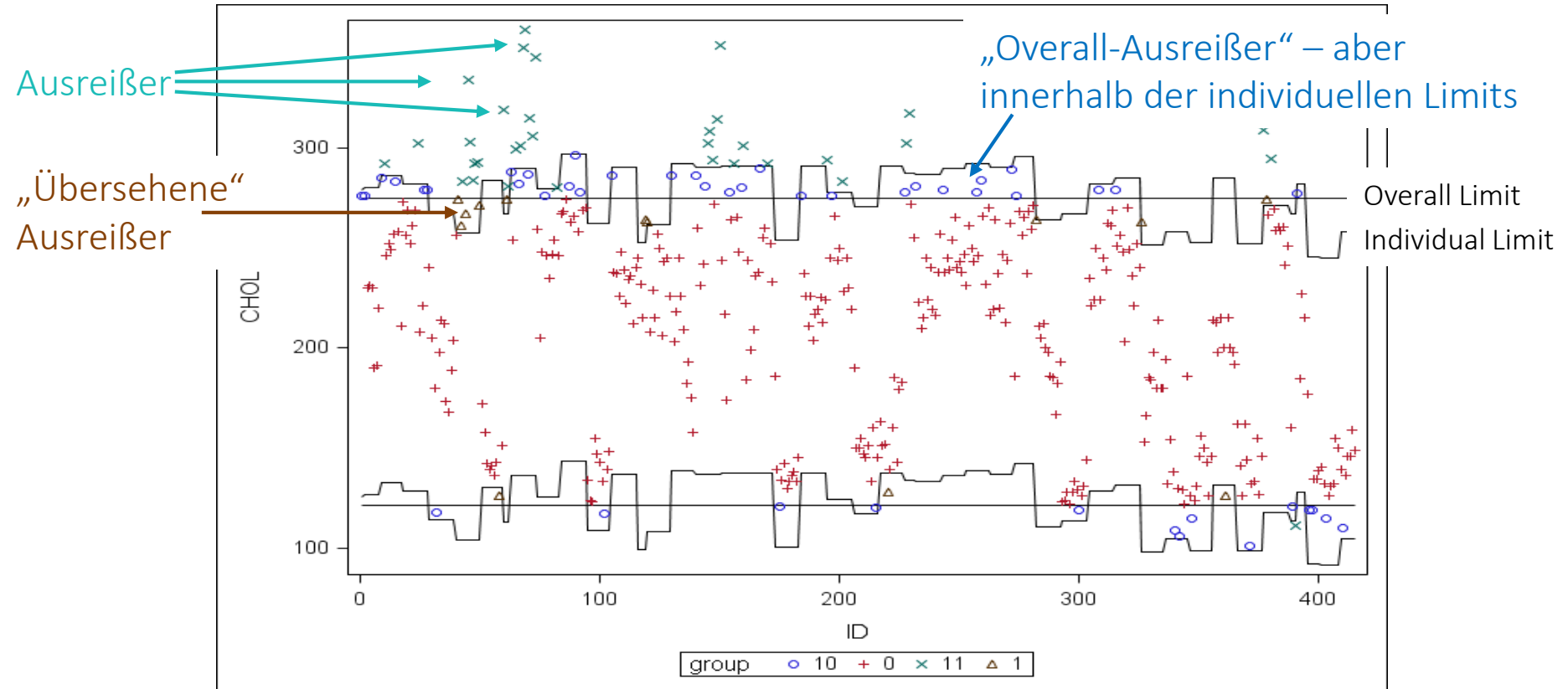


„Ja, aber im Jänner haben wir immer deutlich mehr Ereignisse“

Modell erkennt, dass dieser Wert im Jänner kein Ausreißer ist



„Alle deren Wert größer x ist, sind Ausreißer! - Wirklich?“

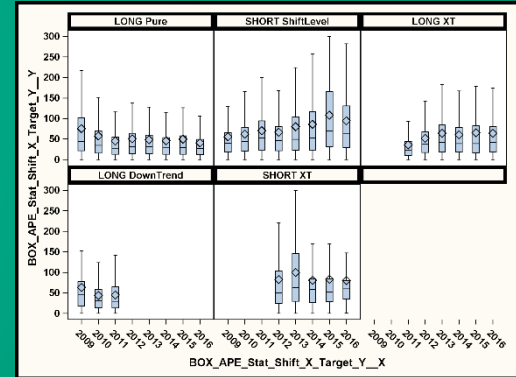
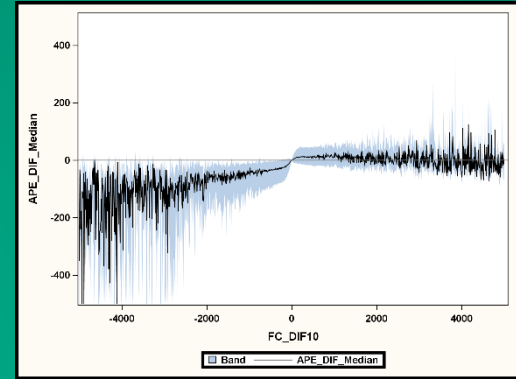


Data Science in Action: #4

Explaining Forecast Errors and Deviations

*Do the demand planners really improve
forecast accuracy with their manual
overwrites?*

Linear Regression
Quantile Regression
Descriptive Statistics

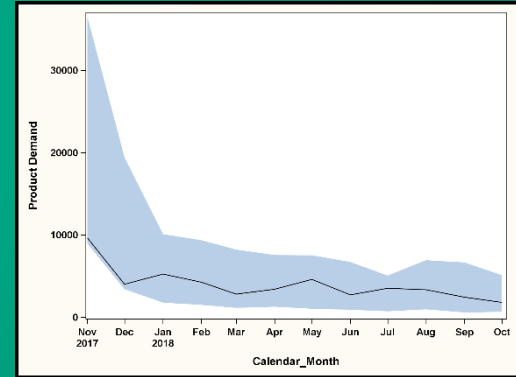
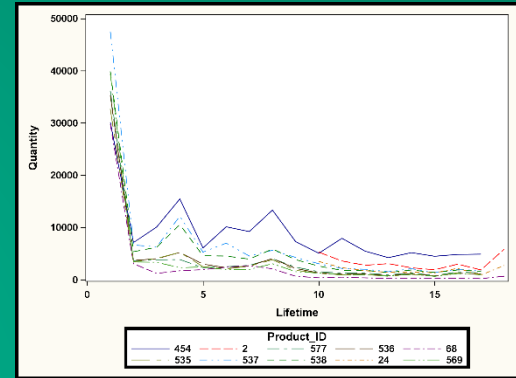


Data Science in Action: #5

Forecasting the Demand for New Products

*Can the expected demand of products
that are introduced only right now be
estimated for forecast planning?*

Poisson Regression
Cluster Analysis
Similarity Search

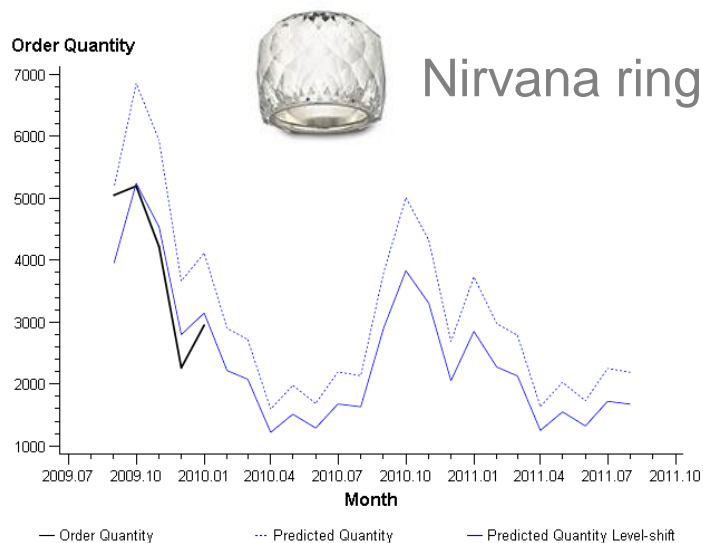


NOVELTY FORECASTING

- Training data from previous collections
- Generalized linear model

- Predictors

- Product attributes
- Time-dependant influence factors
- Number of shops
- Actual order intake
- Actual sell-through



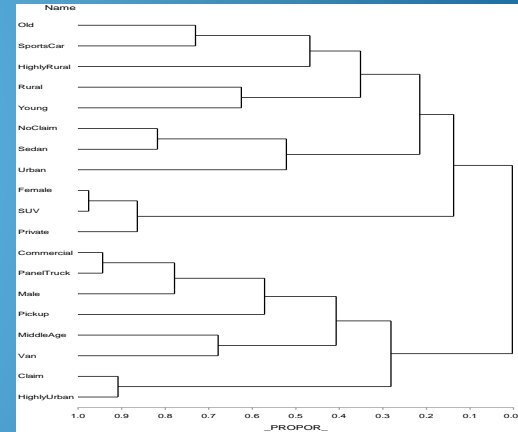
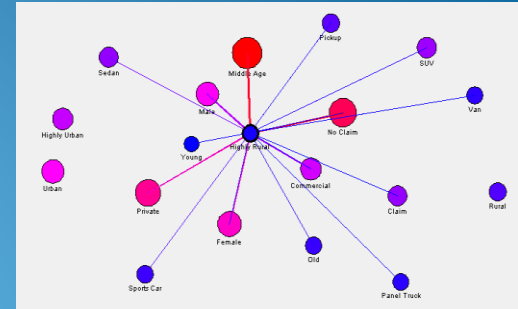
SWAROVSKI

Data Science in Action: #6

Listening to Your Data – Discover Relationships with Unsupervised Analysis Methods

*Can your data tell you stories about
your analysis subjects, even if you don't
ask explicitly?*

Unsupervised machine learning methods:
association analysis
variable clustering



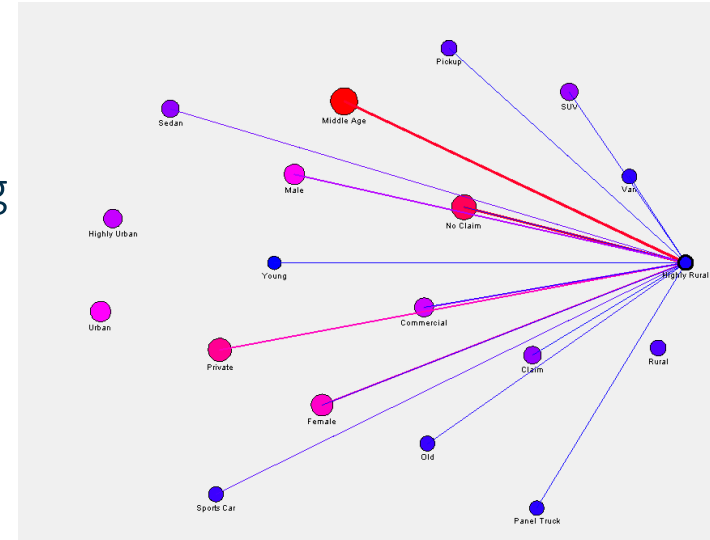
Lassen Sie ihre Daten sprechen!

Auffinden von Zusammenhängen in Ihren Analysedaten

- Daten aus der KFZ-Versicherung mit 6 Eigenschaften pro Versicherungsnehmer

Variable	Feature
AGE	YOUNG, MIDLIFE, OLD
GENDER	MALE, FEMALE
DENSITY	HIGHLY URBAN, URBAN, HIGHLY RURAL, RURAL
CAR_TYPE	VAN, SPORTS CAR, SUV, SEDAN, PICK UP
CAR_USAGE	PRIVATE, COMMERCIAL
CLM_FLAG	CLAIM, NO CLAIM

- Anwenden von unsupervised machine learning (Assoziationsanalyse) um Zusammenhänge zwischen den Eigenschaften aufzudecken.



Trauen Sie sich! Transponieren Sie die Daten, so wie Sie es sonst typischerweise nicht tun.

One-Row-Per-Subject

POLICYNO	CLM_FLAG	CAR_USE	CAR_TYPE	AGE	GENDER	DENSITY
160	No	Private	Sedan	60	M	Highly Urban
24836	No	Commercial	Sedan	43	M	Highly Urban
28046	No	Private	Van	48	M	Urban
28960	No	Private	SUV	35	F	Highly Urban
40933	No	Private	Sedan	51	M	Highly Urban
55277	No	Private	SUV	50	F	Urban
63212	Yes	Commercial	Sports Car	34	F	Highly Urban
69651	No	Private	SUV	54	F	Highly Urban
88070	Yes	Private	Sedan	40	M	Urban
93553	No	Commercial	SUV	44	F	Rural
127444	Yes	Commercial	Van	37	M	Highly Urban
141509	Yes	Private	SUV	34	F	Highly Urban
145326	No	Commercial	Van	50	M	Rural
146809	Yes	Private	Sports Car	53	F	Urban
148250	No	Private	Sedan	43	F	Rural
157851	No	Commercial	Van	55	M	Urban



Multiple-Row-Per-Subject
Key-Value Tabelle

POLICYNO	Feature
160	Highly Urban
160	No Claim
160	Sedan
160	Private
160	Male
160	Old
24836	Highly Urban
24836	No Claim
24836	Sedan
24836	Commercial
24836	Male
24836	Middle Age



Lassen Sie ihre Daten sprechen!

Männer fahren kaum Sportwägen?

Regel 278 enthält, dass Sportwägen nur in 2,54 % der Fälle von Männern gefahren werden (erwartet wären 46 %)



index	RULE	LHAND	RHAND	COUNT	SUPPORT	EXP_CONF	CONF	LIFT
267	Commercial ==> Sports Car	Commercial	Sports Car	200.00	1.94	11.44	5.28	0.46
268	Rural ==> Claim	Rural	Claim	102.00	0.99	26.66	6.52	0.24
269	Claim ==> Rural	Claim	Rural	102.00	0.99	15.18	3.71	0.24
270	Young ==> Highly Urban	Young	Highly Urban	10.00	0.10	34.93	8.33	0.24
271	Highly Rural ==> Claim	Highly Rural	Claim	32.00	0.31	26.66	6.30	0.24
272	Claim ==> Highly Rural	Claim	Highly Rural	32.00	0.31	4.93	1.17	0.24
273	Van ==> Female	Van	Female	117.00	1.14	53.82	12.70	0.24
274	Female ==> Van	Female	Van	117.00	1.14	8.94	2.11	0.24
275	Panel Truck ==> Female	Panel Truck	Female	40.00	0.39	53.82	4.69	0.09
276	Male ==> SUV	Male	SUV	99.00	0.96	27.98	2.08	0.07
277	SUV ==> Male	SUV	Male	99.00	0.96	46.18	3.43	0.07
278	Sports Car ==> Male	Sports Car	Male	30.00	0.29	46.18	2.54	0.06

- Kann anzeigen, dass in unserer Datenbasis tatsächlich Sportwägen in erster Linie von Frauen gefahren worden.
- Kann auch ein Trigger für eine detailliertere Analyse der Datenqualität sein.
- Ein fachliche Erklärung kann sein, dass der Sportwagen das 2. oder 3. Auto in der Familie ist, und dieser aus steuerlichen Gründe auf die Ehefrau registriert ist.
- Möglicherweise bietet ein Mitbewerber eine Polizze für Männer zu einem deutlich besseren Preis an.

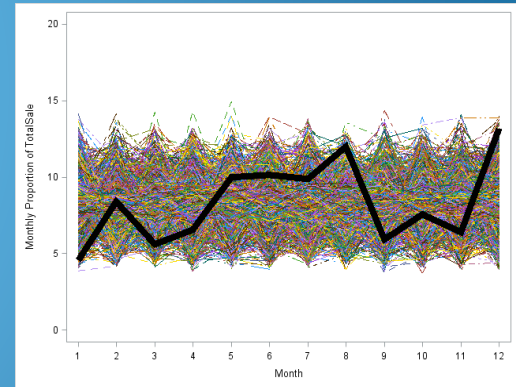
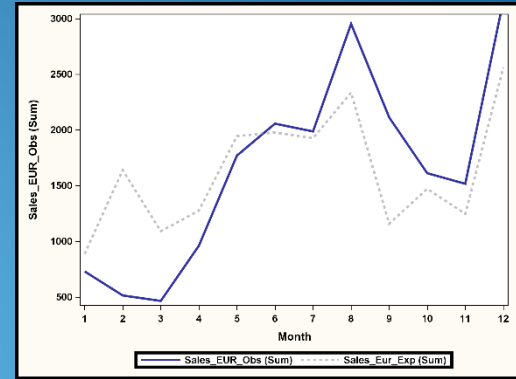


Data Science in Action: #7

Checking the Alignment with Predefined Pattern

*Which customers show a behavior that
is far from what you expected?*

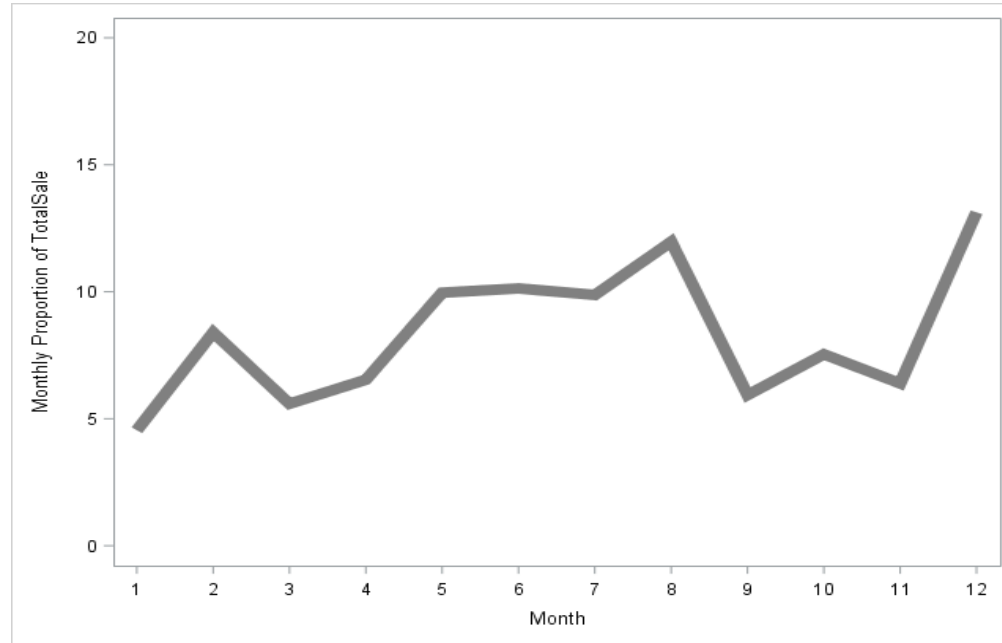
Chi2 independency test
Benford's law
Time Series Similarity



“Welche meiner Verkäufer halten sich kaum an unsere Vorgaben?”

Der Bedarf an “Sub-Contracts” für ein Cateringunternehmen variiert im Verlauf eines Kalenderjahres

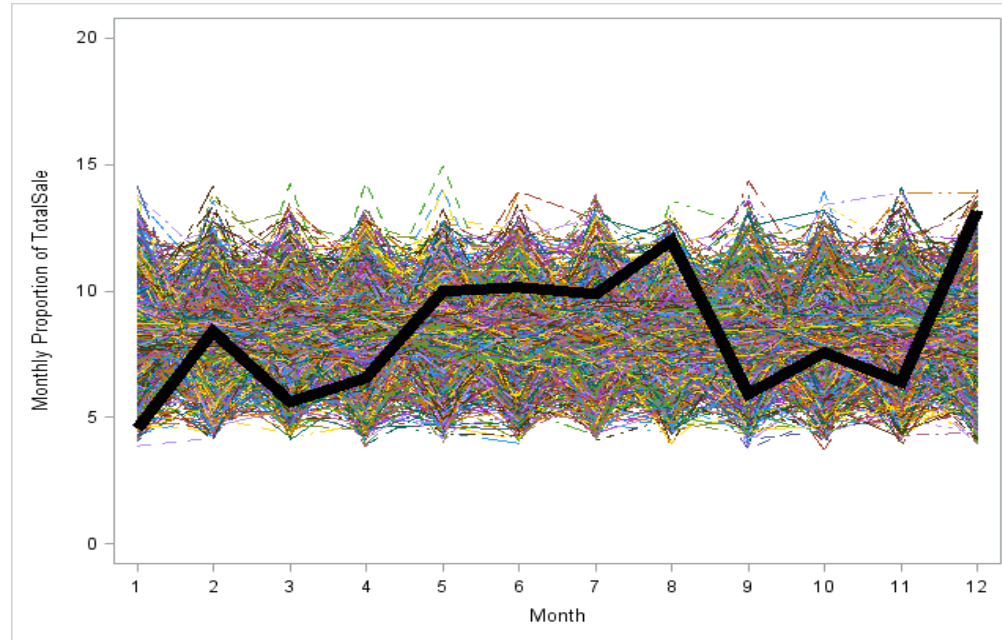
Verkäufer sind angehalten, entsprechend dieses Musters Verträge zu akquirieren.



Anzeige der Jahresverläufe pro Verkäufer hilft nicht wirklich

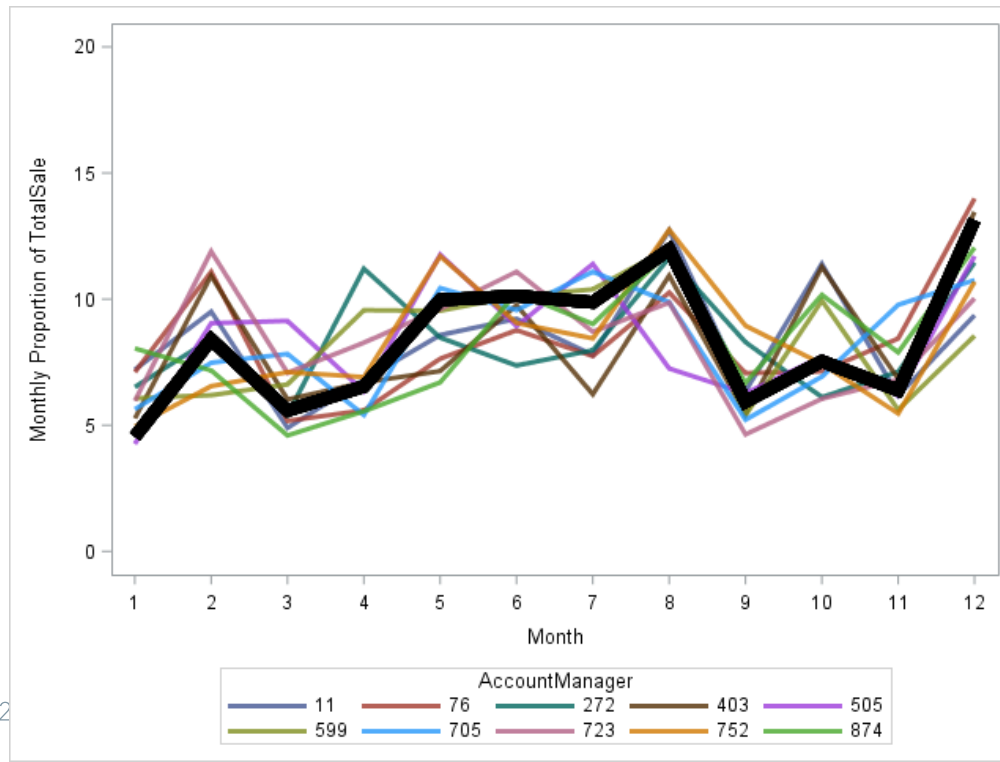
Kein klares Bild.

Unmöglich,
alle Linien einzeln
durchzusehen.



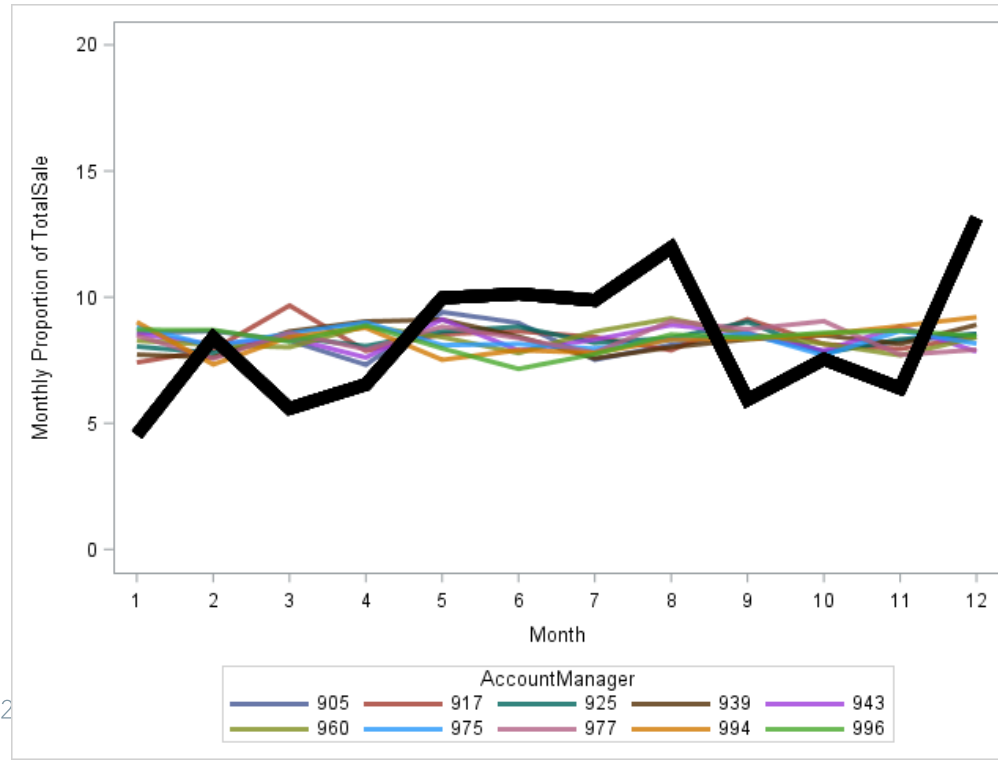
Ranking der Verkäufer mit analytischen Methoden (1)

Top 10 Verkäufer bzgl. "Alignment" mit der Vorgabe



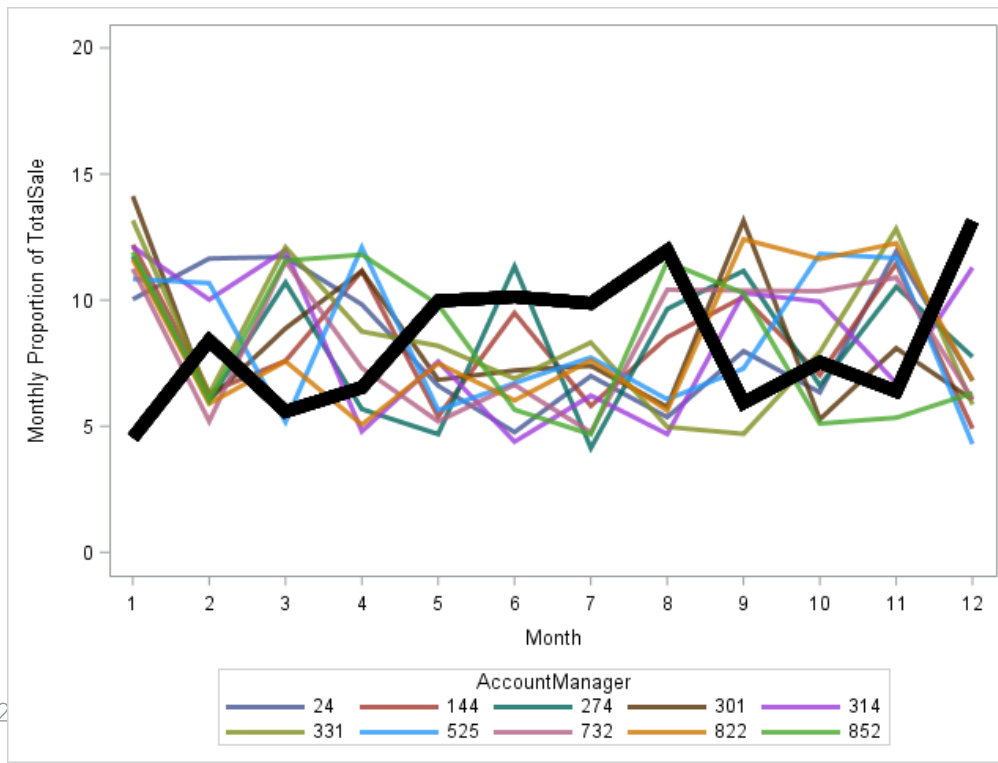
Ranking der Verkäufer mit analytischen Methoden (2)

Top 10 Verkäufer, für die es keine saisonale Variation gibt.

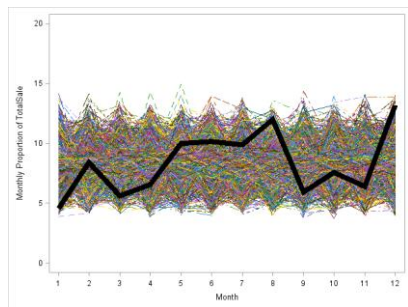


Ranking der Verkäufer mit analytischen Methoden (3)

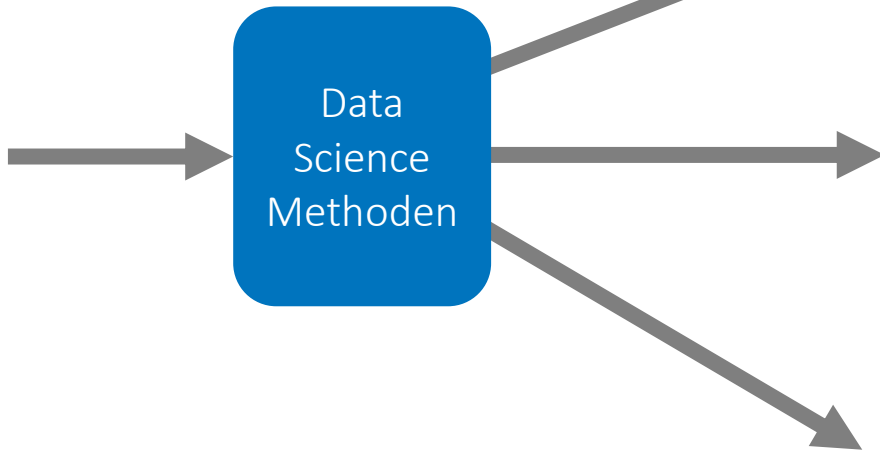
Top 10 Verkäufer die “gegen” das Muster arbeiten



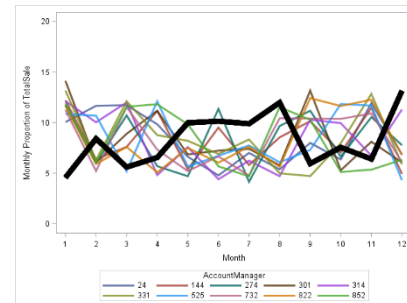
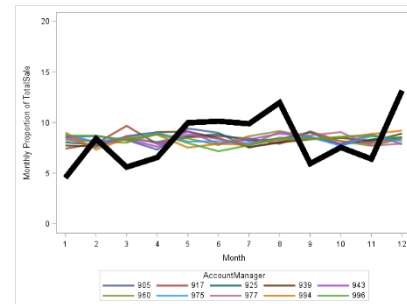
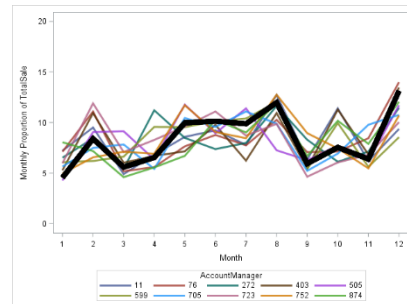
Analytik hilft mir, ein klareres Bild zu gewinnen!



Vom „Rauschen“



zu interpretierbaren
Segmenten



Topic Search Documents and Clustering

- Text Mining
- Text Parsing (Synonyme, Stemming, Stop-Listen)
- Term by Document Weights



Kann ich ähnliche Kapitel erkennen, ohne die Bücher (von Gerhard 😊) erst lesen zu müssen?

Topic > +access,+file,+text,+relational,+relational database



PAGE 104 **Data Preparation** for Analytics Using SAS Chapter 13: **Accessing Data** PAGE 103 Part 3 **Data Mart Coding** and Content Chapter 13 **Accessing Data** Transposing One- and Multiple-Rows-per-Subject **Data Structures** 115 Chapter 15 Transposing **Longitudinal Data** 131 Chapter 16 **Transformations of** Chapter 17 **Transformations of Categorical Variables** 161 Chapter 18 Multiple **Interval-Scaled** Observations per **Subject** 179 Chapter 19 **Multiple Catego**



PAGE 38 **Data Preparation** for Analytics Using SAS Chapter 5: The **Origin of Data** PAGE 43 Part 2 **Data Structures** and **Data Modeling** Chapter 5 The **Models** 45 Chapter 7 **Analysis Subjects** and Multiple Observations 51 Chapter 8 The One-Row-per-Subject **Data Mart** 61 Chapter 9 The Multiple-Rows-p **Data Structures** for **Longitudinal** Analysis 77 Chapter 11 Considerations for **Data Marts** 89 Chapter 12 Considerations for Predictive **Modeling** 95 Introdu



PAGE 178 **Data Preparation** for Analytics Using SAS Chapter 17: **Transformations of Categorical** Variables PAGE 177 Chapter 17 **Transformations** Introduction 17.2 General Considerations for **Categorical** Variables 162 17.3 **Derived** Variables 164 17.4 Combining **Categories** 166 17.5 **Dummy Coding** **Multidimensional Categorical** Variables 172 17.7 **Lookup Tables** and **External Data** 176 17.1 Introduction In this chapter we will deal with **transformatio**



40 **Data** Quality for Analytics Using SAS Chapter 3: **Data** Availability 41 Chapter 3: **Data** Availability 3.1 Introduction 32 3.2 General Considerations 32 Re: **data** availability 32 Availability and usability 32 Effort to make **data** available 33 Dependence on the **operational process** 33 Availability and alignment in t of Historic **Data** 34 **Categorization** and examples of historic **data** 34 The **length** of the **history** 35 **Customer event histories** 35 **Operational systems** and a



PAGE 382 **Data Preparation** for Analytics Using SAS Appendix B: The Power of **SAS** for **Analytic Data Preparation** PAGE 381 Appendix B The Power c 369B.1 Motivation B.2 Overview 370 B.3 **Extracting Data** from **Source Systems** 371 B.4 Changing the **Data Mart Structure**: Transposing 371 B.5 **Data Mar** Multiple-Rows-per-Subject **Data** Sets 372 B.6 Selected Features of the **SAS Language** for **Data Management** 375 B.7 Benefits of the **SAS Macro Langu**



#SASF17



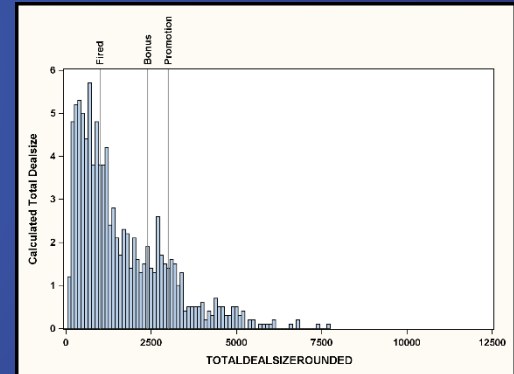
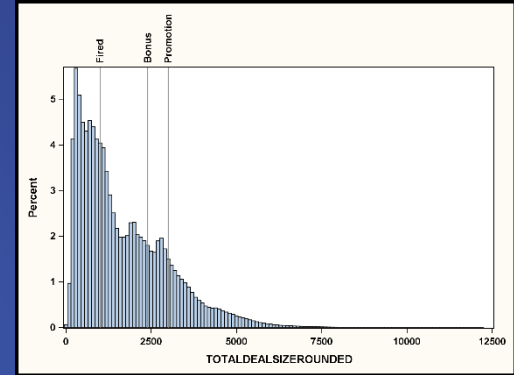
SAS

Data Science in Action: #9

Using Monte Carlo Simulations to Understand the Outcome Distribution

*When the sales manager looks at the
project pipeline, does the sum of weighted
averages give him or her a full picture?*

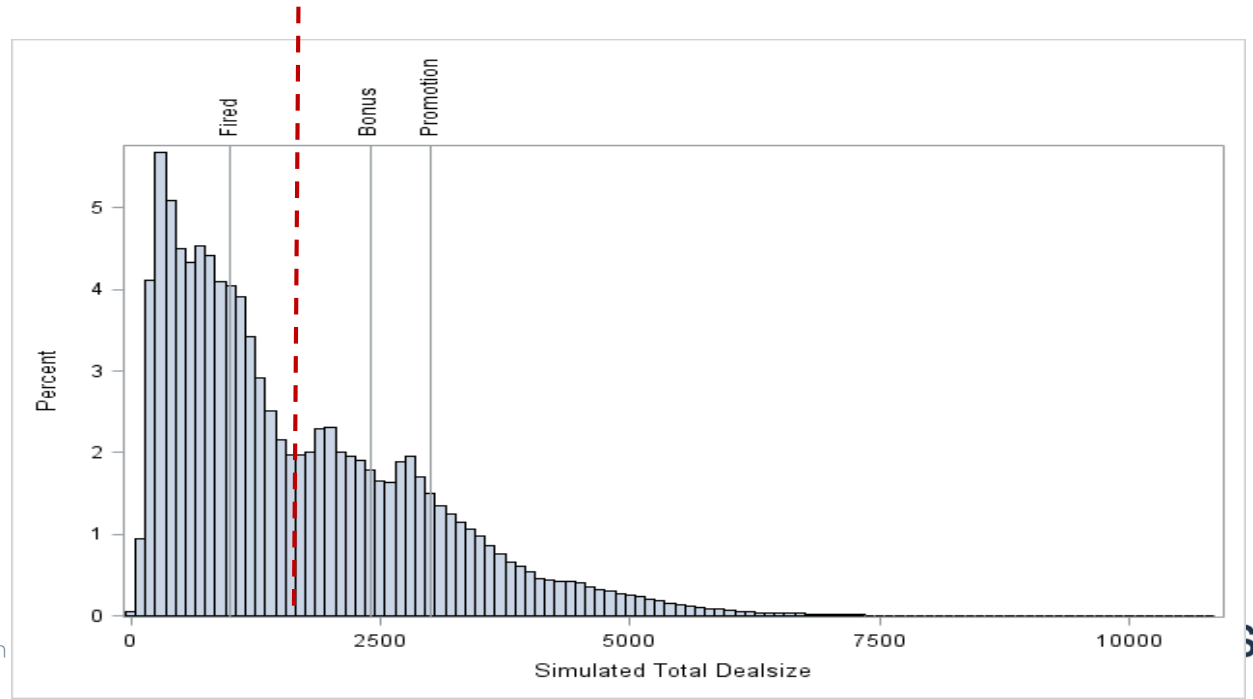
Monte Carlo simulations
Mathematical programming



Wird der Sales Manager seinen Job behalten?

ProjectID	DealSize (1000 \$)	Proba- bility
1	1500	10%
2	10	65%
3	500	20%
4	50	50%
5	100	40%
6	30	90%
7	10	60%
8	150	20%
9	200	25%
10	180	10%
11	900	10%
12	750	20%
13	600	10%
14	320	20%
15	100	40%
16	50	80%
17	2000	5%
18	400	20%
19	2500	10%
20	1700	15%
21	100	80%

Gewichtetes Mittel:
\$ 1.661.500

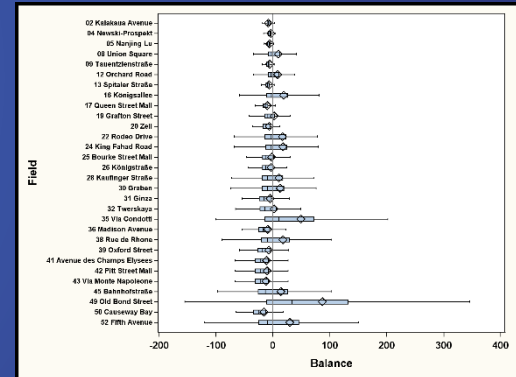
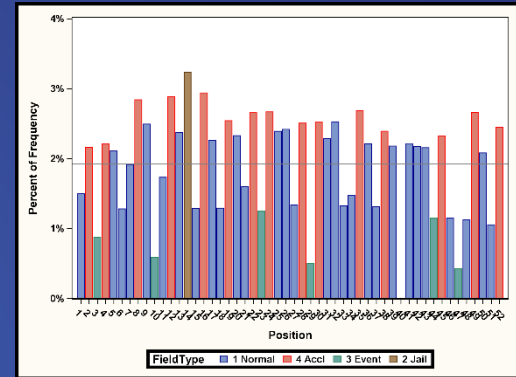


Data Science in Action: #10

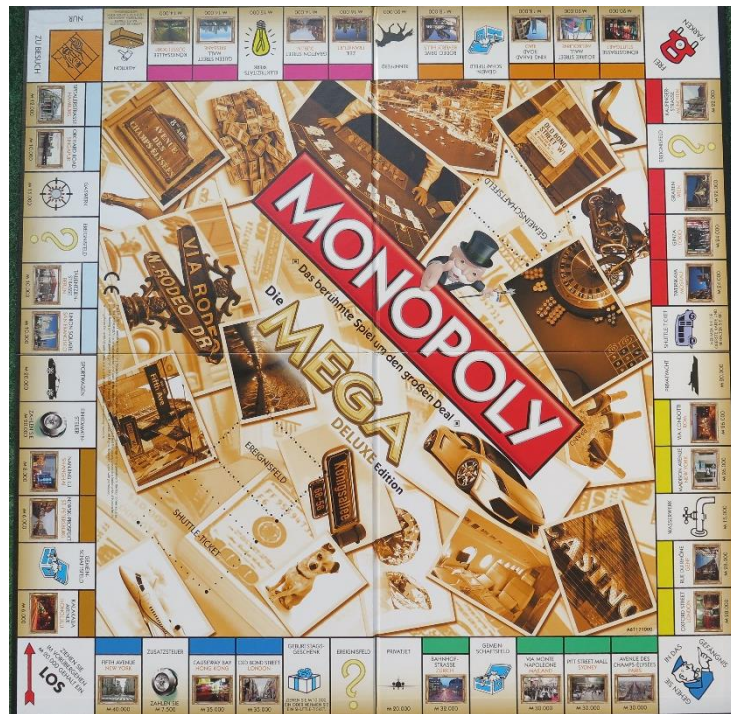
Studying Complex Systems – Simulating the Monopoly Board Game

*How can you simulate complex
environments to get insight in the most
frequent processes?*

Monte Carlo Simulations



Das Monopoly Spiel ist vielen Frameworks im Geschäftsleben gar nicht so unähnlich



Monetäre Dimension



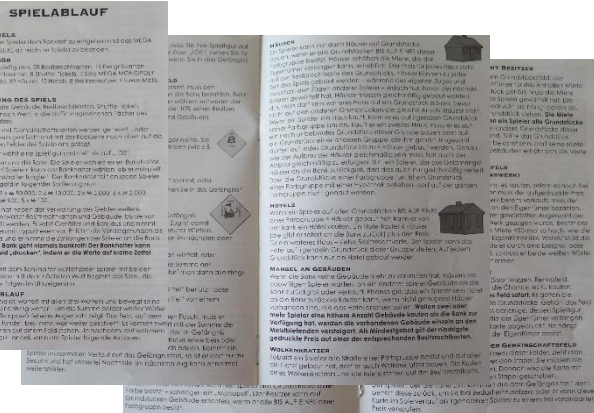
Dynamische Komponenten

Zufällige

Komponenten



Rahmenwerk von Möglichkeiten und Ereignissen



Komplexe Regeln

Zusätzliche Anweisungen

Simulation komplexer Prozesse erlaubt mir Einblick in Zusammenhänge (die ich sonst nicht gesehen hätte)



Würfel-
Summe



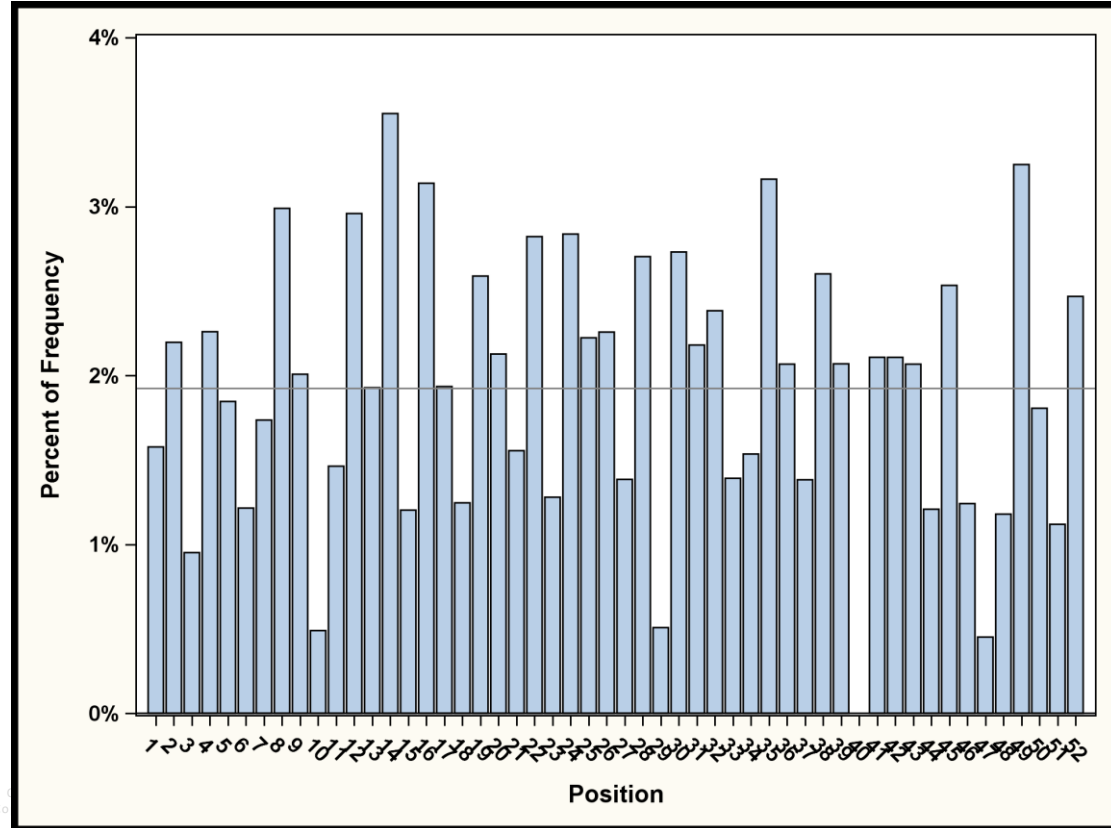
Gehe ins
Gefängnis!



Ereignis
Felder

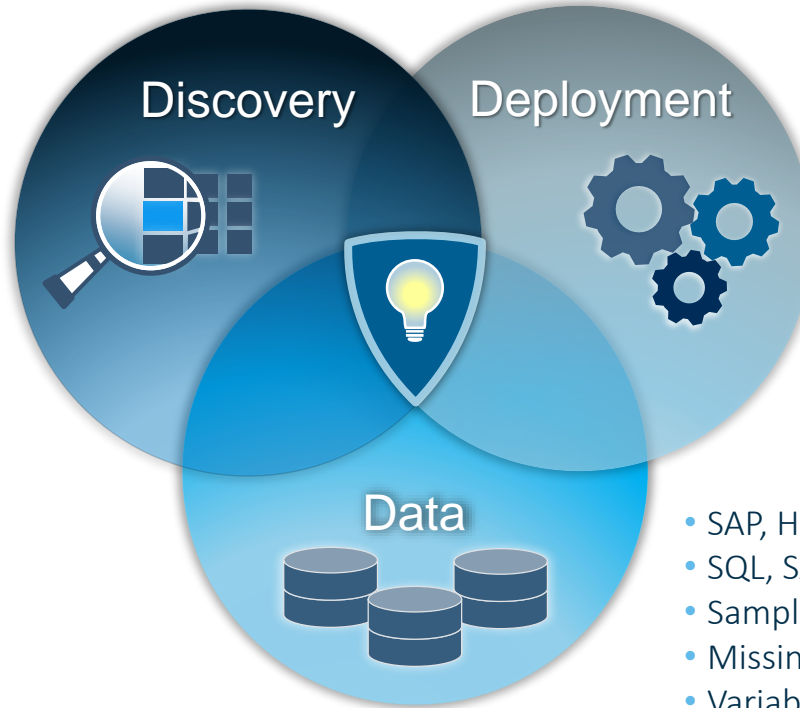


Accelerator
Würfel



Data Mining und Machine Learning mit der SAS Analytic Plattform

- Logistic Regression
- Linear Regression
- Generalized Linear Models
- Nonlinear Regression
- Ordinary Least Squares Regression
- Decision Trees
- Partial Least Squares Regression
- Quantile Regression
- K-means and K-modes Clustering
- Principal Component Analysis
- Random Forest
- Gradient Boosting
- Neural Networks
- Support Vector Machines
- Factorization Machines
- Network Analytics/Community Detection
- Text Mining
- Boolean Rules
- Auto-tuned Hyper-parameters



- Assess Supervised Models
- Modellverwaltung
- Deployment
- Laufende Validierung
- Modell-Retirement
- Retraining

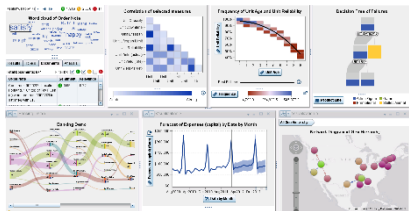
- SAP, Hadoop, Streaming, rel.DB, ...
- SQL, SAS Datastep, Matrix
- Sampling and Partitioning
- Missing Value Imputation
- Variable Binning
- Variable Selection
- Transpose

Offenheit der SAS Analytic Plattform für unterschiedliche Zugriffsarten

Erfüllung der individuellen Anforderungen



Office
Integration



One Integrated Solution for Different User Types



#SASF17



SAS FORUM 2017 / sasforum17

Copyright © SAS Institute Inc. All rights reserved.



Key Takeaways

Analytics und Data Science sind da um Ihnen zu helfen!

- Sie sehen ein klareres, objektiveres Bild Ihrer Daten und Analyse-Subjekte
- Sie erhalten explizite Ergebnisse anstatt die Nadel im Heuhaufen zu suchen
- Die Daten sprechen zu Ihnen und Sie erhalten die Ergebnisse automatisch statt manuell
- Do it again! – Behandeln Sie Ihre Modelle als “Asset” und wiederholen Sie Ihre Analyse

Machine Learning and Data Science sind das Kernstück der SAS Analytic Platform

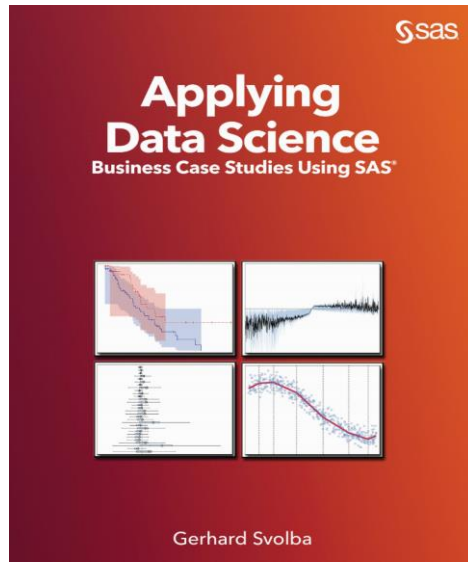
- Umfassendes Set an Methoden – Entdecken und Produktivstellen
- Offen für unterschiedliche Benutzertypen (Coding, Point&Click, SAS, R, Python, ...)



More Information

Gerhard Svolba – Principal Analytic Solutions Architect

sastools.by.gerhard@gmx.net



- Applying Data Science – Business Case Studies Using SAS, SAS Press 2017
- Eight Case Studies showing how Data Science and Analytics can be applied to provide insight into your data and improve your business decisions
- [http://www.sascommunity.org/wiki/Applying_Data_Science - Business Case Studies Using SAS](http://www.sascommunity.org/wiki/Applying_Data_Science_-_Business_Case_Studies_Using_SAS)



#SASF17



SAS F

Aktuellste Version des Vortrags → Google → Gerhard sas samples



Further Links

- Gerhard Svolba: Mehr als linear oder logistisch – ausgewählte Möglichkeiten neuer Regressionsmethoden in SAS - Download the [presentation](#) and the [paper](#)
- Allison, P. 1995. *Survival Analysis Using SAS®: A Practical Guide, Second Edition*. Cary, NC: SAS Institute Inc.
- SAS/STAT® 14.2 User's Guide. The LIFETEST Procedure.
<http://support.sas.com/documentation/onlinedoc/stat/142/lifetest.pdf>
(accessed 1 March 2017).
- Kuhfeld, W., and W. Cai. 2013. "Introducing the New ADAPTIVEREG Procedure for Adaptive Regression." SAS Global Forum Proceedings.
<http://support.sas.com/resources/papers/proceedings13/457-2013.pdf>
(Paper 457-2013).

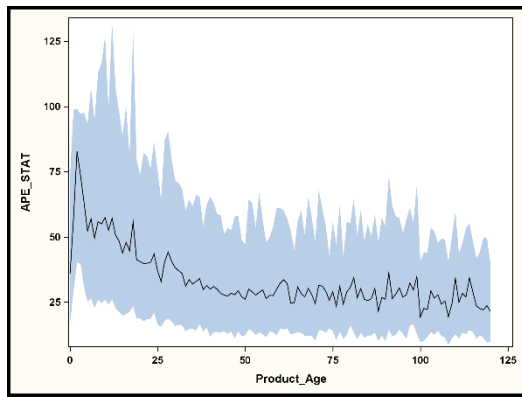


An welchen Schrauben soll ich drehen, wenn ich die Modellqualität verbessern will?

Univariates Modell

Ranking	Input Variable	R Square linear
1	MODEL	0.0554
2	PRODUCT_AGE	0.0433
3	PRODUCT_GROUP	0.0224
4	LAUNCH_MONTH	0.0172
5	TARGET_YEAR	0.0102
6	TARGET_CALMONTH	0.0084
7	LEAD_TIME	0.0046
8	PRICE_INDEX	0.0016

Stepwise Selection Summary					
Step	Effect Entered	Effect Removed	Number Effects In	Number Parms In	Adjusted R-Square
0	Intercept		1	1	0.0000
1	Model		2	5	0.0533
2	Product_Group		3	18	0.0611
3	Target_CalMonth		4	29	0.0682
4	Product_Age		5	30	0.0748
5	Lead_Time		6	31	0.0831
6	target_year_shift		7	32	0.0856
7	Launch_Month		8	43	0.0865
8	Price_Index		9	44	0.0869*



Variable (Category)			Coefficient
Intercept			48.730894
Model	LONG Pure		-11.806946
Model	SHORT XT		17.271517
Model	LONG DownTrend		-5.715497
Model	LONG XT		-10.407650
Model	SHORT ShiftLevel		10.6586
Target_CalMonth	1		-13.700216
Target_CalMonth	2		-1.941132
Target_CalMonth	3		1.000407

